

Abstract

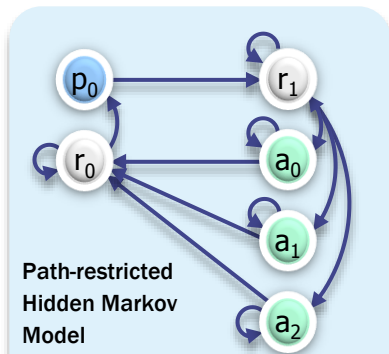
Linear Predictive Coding (LPC), proposed in the 60s, has established the modern speech analysis/synthesis framework and has opened the door of today's mobile and VoIP communication technology. While LPC has realized the analysis/synthesis framework focusing on the 'phonemic' factor of speech, the aim of this work is to develop a new analysis/synthesis framework focusing on the 'prosodic' factor. Although a well-founded physical model for vocal fold vibration was proposed in the 60s by Fujisaki (known as the "Fujisaki model"), how to estimate the underlying parameters has long been a difficult task. We have developed a stochastic counterpart of the Fujisaki model, which made it possible to apply powerful statistical inference techniques to accurately estimate the underlying parameters. This model has a high potential to be developed into a next-generation module for Text-to-Speech, speech analysis, synthesis and conversion systems.

What is fundamental frequency (FO) contour?

- An acoustic correlate that plays an important role in conveying non-linguistic information such as the identity, intention, attitude and mood of the speaker

The aim of this work

- Develop stochastic counterpart of the Fujisaki model
→ Elegant parameter estimation framework
- Realize an analysis and synthesis framework of speech prosody and natural-sounding Text-to-Speech synthesis



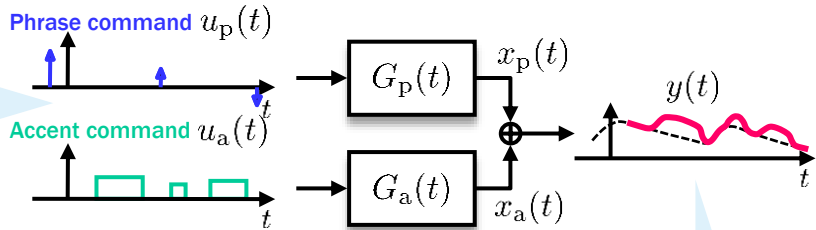
For $t = 1, \dots, T$:

$s_t | s_{t-1} \sim \pi_{s_{t-1}, s_t}$ (state sequence)

$$\begin{bmatrix} u_p[t] \\ u_a[t] \end{bmatrix} | s_t \sim \mathcal{N} \left(\begin{bmatrix} \mu_{p, s_t} \\ \mu_{a, s_t} \end{bmatrix}, \begin{bmatrix} \sigma_p^2 & 0 \\ 0 & \sigma_a^2 \end{bmatrix} \right)$$

Change of variables

Physical model for FO generation process (Fujisaki model)

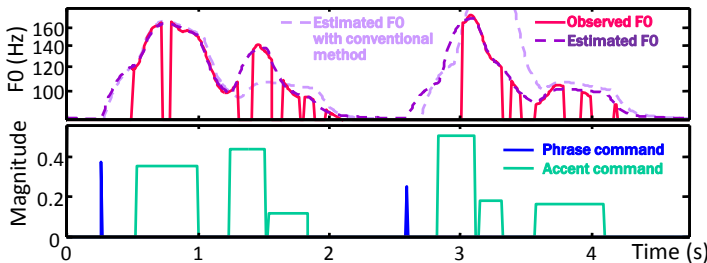


Generative model of FO contours

$$y \sim \mathcal{N}(\mu, \Sigma)$$

$$\begin{cases} \mu = G_p \mu_p + G_a \mu_a + \mu_b \mathbf{1} \\ \Sigma = \sigma_p^2 G_p G_p^T + \sigma_a^2 G_a G_a^T + \Sigma_b \end{cases}$$

Estimated phrase and accent commands with proposed method



What will become possible?

- Manually changing phrase and accent commands allows to convert intonation as desired while keeping naturalness of speech
- Predicting phrase and accent commands from text inputs allows to synthesize speech with natural-sounding FO contours

Audio demo

Related works

- [1] H. Kameoka, K. Yoshizato, T. Ishihara, K. Kadowaki, Y. Ohishi, and K. Kashino, "Generative modeling of voice fundamental frequency contours," *IEEE/ACM Transactions on Audio, Speech and Language Processing*, to appear, 2015.
- [2] K. Kadowaki, T. Ishihara, N. Hojo, and H. Kameoka, "Speech prosody generation for text-to-speech synthesis based on generative model of FO contours," in *Proc. The 15th Annual Conference of the International Speech Communication Association (Interspeech 2014)*, pp. 2322-2326, Sep. 2014.

Contact

Hirokazu Kameoka Recognition Research Group, Media Information Laboratory
E-mail : kameoka.hirokazu(at)lab.ntt.co.jp

